



## Report on the Code

### PD-Internship

Lerner Paul

25/04/2018

Some of this report was originally present in report #4. However, since it's prone to update unlike the rest of report #4, on April 19th I split the report in 2 : Literature review and Code. This is the Code one.

In parallel I read about online HaW analysis and thus wrote report #5.

Disambiguation : All the results were obtained on the *spiral* task if not specified otherwise.

All the experimental results were obtained while validating on the test set. See [5 Evaluation and experiments](#) for discussion. Cf. Git Commits and the experiment results.

# Summary

<b>1 Model architecture</b>	<b>1</b>
<b>2 Model training</b>	<b>6</b>
<b>3 Regularization</b>	<b>7</b>
<b>4 Data preprocessing and representation</b>	<b>9</b>
<b>5 Evaluation and experiments</b>	<b>11</b>
<b>6 Visualization &amp; Interpretation</b>	<b>12</b>
6.1 Input weights	12
6.2 Forget Gate	13
6.3 Misclassified subjects	14
<b>Conclusion - Todo List</b>	<b>14</b>

# 1 Model architecture

We use the LSTM as defined in PyTorch<sup>1</sup>, thus, for each element in the input sequence, each LSTM unit of each layer computes the following operations :

$$\begin{aligned} i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{(t-1)} + b_{hi}) \\ f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{(t-1)} + b_{hf}) \\ g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{(t-1)} + b_{hg}) \\ o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{(t-1)} + b_{ho}) \\ c_t &= f_t c_{(t-1)} + i_t g_t \\ h_t &= o_t \tanh(c_t) \end{aligned}$$

Where  $h_t$  is the hidden state at time  $t$ ,  $c_t$  is the cell state at time  $t$ ,  $x_t$  is the input at time  $t$ ,  $h_{(t-1)}$  is the hidden state of the layer at time  $t-1$ , and  $i_t, f_t, g_t, o_t$  are the input, forget, cell, and output gates, respectively.  $\sigma$  is the sigmoid function.  $W_*$  are the weight matrices and  $b_*$  are the bias vectors.

**Baseline decoder :** We currently feed the last hidden state of the LSTM (summed if bLSTM) to a FC layer with a single neuron which is used for binary classification after applying a Sigmoid to it.

Our dataset is challenging :

1. by its very small number of sample
2. by its very long and varying sequences

Thus our model should be complex enough to capture the long sequences but not too much or heavily regularized so it won't overfit the very little number of samples.

With a very simple model (**baseline**) :

- *learning\_rate* =  $10^{-3}$
- *hidden\_size* = 100
- *num\_layers* = 1
- *bidirectional* = False
- *dropout* = 0.0
- *no gradient clipping*
- → **43 701** total parameters, all trainable.

We're able to fit the training set (90% of the dataset) in 30 epochs until *loss* = 0.323, *accuracy* = 0.906

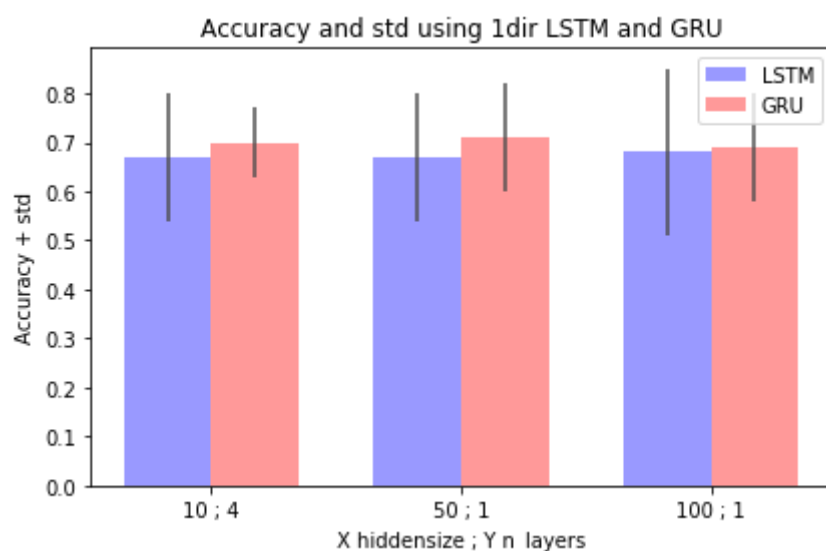
---

<sup>1</sup> <https://pytorch.org/docs/stable/nn.html#torch.nn.LSTM>

Although the model starts overfitting after 1, 2, 1, 4, 4, 2, 4, 3, 5 and 1 epochs (out of 10 folds, respectively).

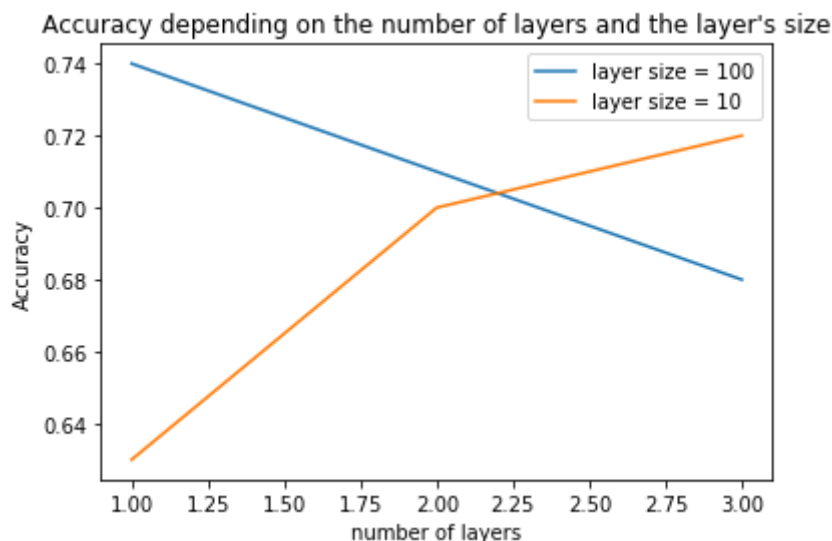
Using a bidirectional LSTM consistently improved the results over the *spiral* task.

Using a 1d GRU consistently improved results and decreased std over 1d LSTM :



There is a lot of interaction between the layer size and the number of layers, we can observe a reversed effect, see figure below. The other hyperparameters are fixed to :

dowsampling factor	learning_rate	bidirectional	carry over	dropout	gradient clipping
1	0.001	TRUE	0	0.5	5

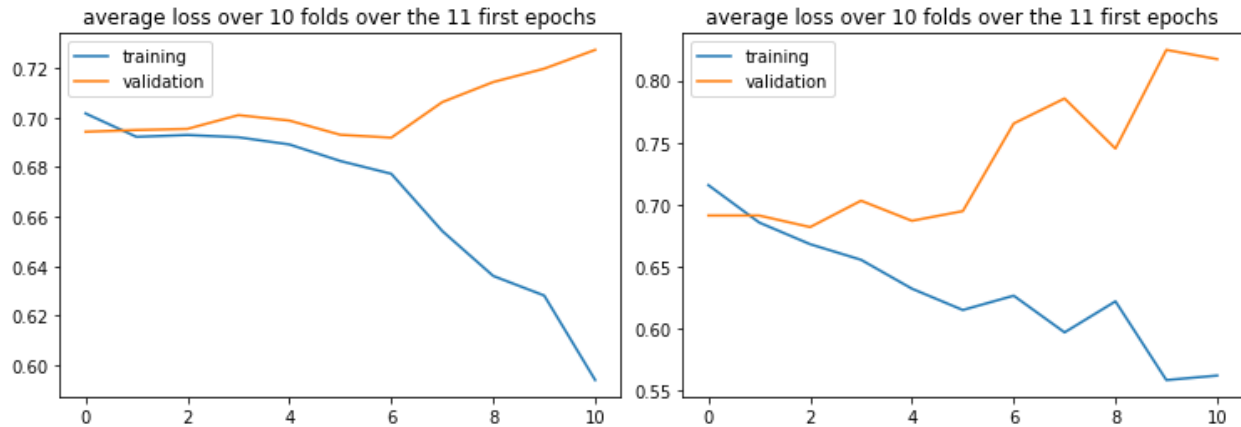


This can probably be explained by the n° of parameters :

layer's size	n° layers	n° of parameters
10	1	1531
10	2	4091
<b>10</b>	<b>3</b>	<b>6651</b>
<b>100</b>	<b>1</b>	<b>87301</b>
100	2	328901
100	3	570501

Therefore, the optimal n° of parameters of the network is probably between 6k and 90k, when using the current regularization techniques (i.e. 50% feedforward dropout). To confirm this hypothesis, I tried with *hidden\_size=50 num\_layers=1* → **23651** parameters and achieved similar results. However with *hidden\_size=50 num\_layers=2* → **84451** parameters the results were slightly inferior. Thus we should privilege the number of layers over the hidden size for the same numbers of parameters. Although this could be explained because of our regularization techniques (i.e. 50% feedforward dropout), if we'd apply recurrent dropout the results might be different.

Among this similar results, I prefer to stick with the simpler (i.e. with less parameters) model which is less prone to overfitting :

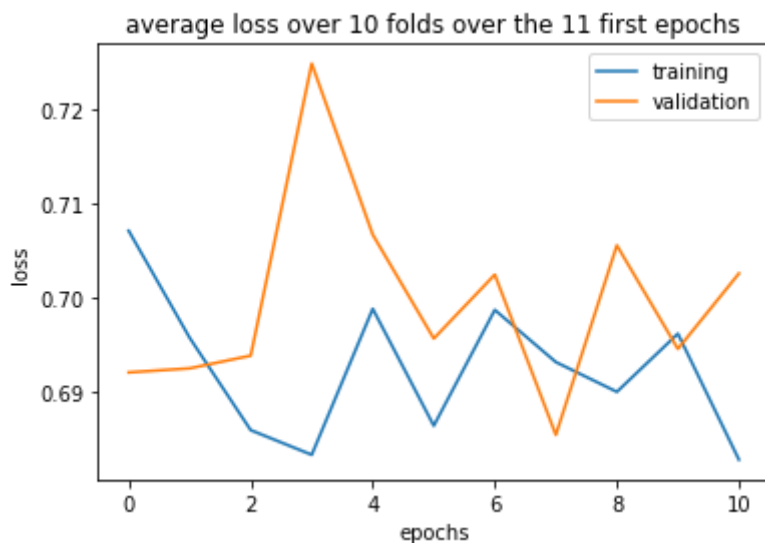


**Fig.** Left :  $hidden\_size=10$   $num\_layers=3$ , right :  $hidden\_size=100$   $num\_layers=1$ . Different scales

I performed a random search on the  $l$  task on :

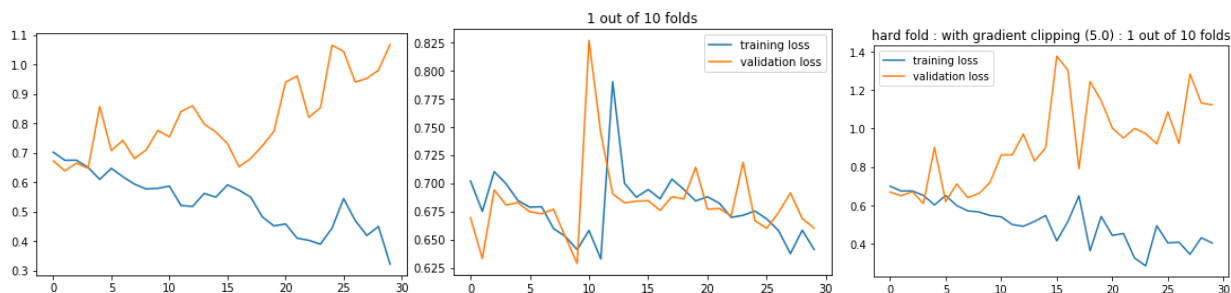
- $learning\_rate$  in  $\{0.01, 0.001, 0.0001\}$
- $num\_layers$  in  $\{1, 2, 3, 4, 5\}$
- $hidden\_size$  in  $[[1, 200]]$
- $bidirectional$  in  $\{True, False\}$
- $dropout$  in  $[0, 1]$

With an extra constraint :  $num\_layers * hidden\_size < 300$  in order to limit the n° of parameters. The random search performed 59 experiments (10-CV) in 3 days. Unfortunately it didn't get better than my manual fine-tuning. We should be very careful of the cross validation results because of the early stopping. One set of parameters achieved 70% validation accuracy but had only 51% training accuracy so I had a look and the model was almost random (cf figure below). Almost all experiments results have better training than validation accuracy. Although this is partly explained by the early stopping (cf. [5 Evaluation and experiments](#)), it is also because we first train the model before validating so a more relevant comparison would be between validation accuracy at epoch  $e$  and training accuracy at epoch  $e + 1$ .



**Fig. Model that doesn't fit but get's 70% accuracy**

## 2 Model training

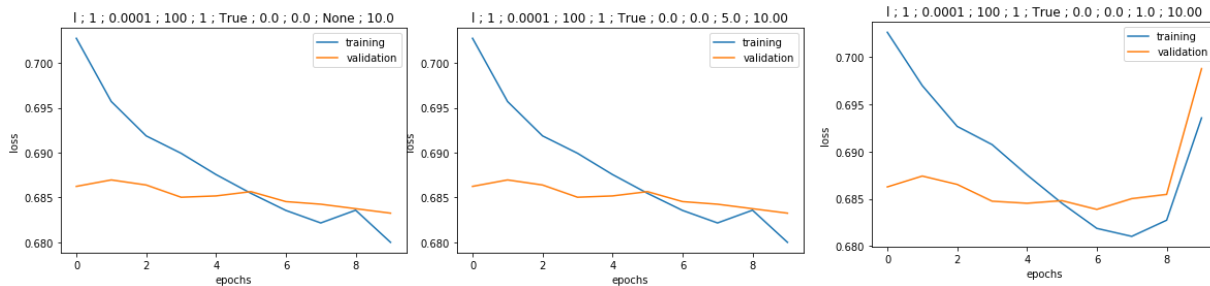


**Fig. Baseline model (1 Model architecture) : on some folds the model won't fit the data very well (left : typical fold, middle : hard fold, right : clipped hard fold) ! \ Different scales ! \**

**It may be because of exploding gradients** : clipping the gradients norm to 5 helped with this (cf. figure).

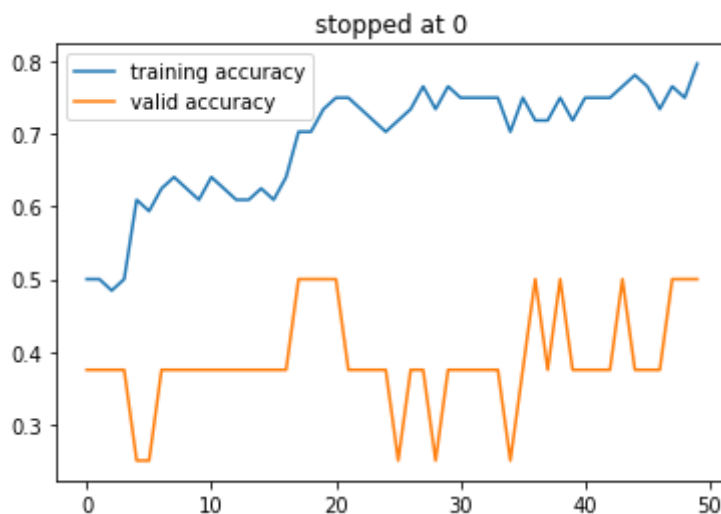
Although it didn't help with the cross validation results, on the contrary, they are slightly inferior compared to the ones without gradient clipping. The differences between these results might also be explained by the random initialization of the weights of the network. I set the random seed after the first few experiments to avoid this bias.

Unlike with *spiral* above, clipping the gradients to 5 didn't have any effect on the *l* task (with  $lr = 10^{-3}$  or  $10^{-4}$ ).



**Fig. Left : no gradient clipping, middle : gradient clipping at 5, right : gradient clipping at 1.**

According to these first experiments, the best learning rate is  $10^{-3}$ . With  $10^{-4}$  there is more variance (std) and on some folds the model is overfitting from the first epoch (cf. figure below), we might try it again along with regularization techniques.



**Fig. LR =  $10^{-4}$ , model is overfitting from the first epoch**

I tried to train the model "continuously" by "passing on" the hidden state from one subject to another (i.e. the hidden state is not init to 0 but to the previous hidden state).

This makes the loss more unstable as you can see below. I highlighted the results in **purple** in the results summary.

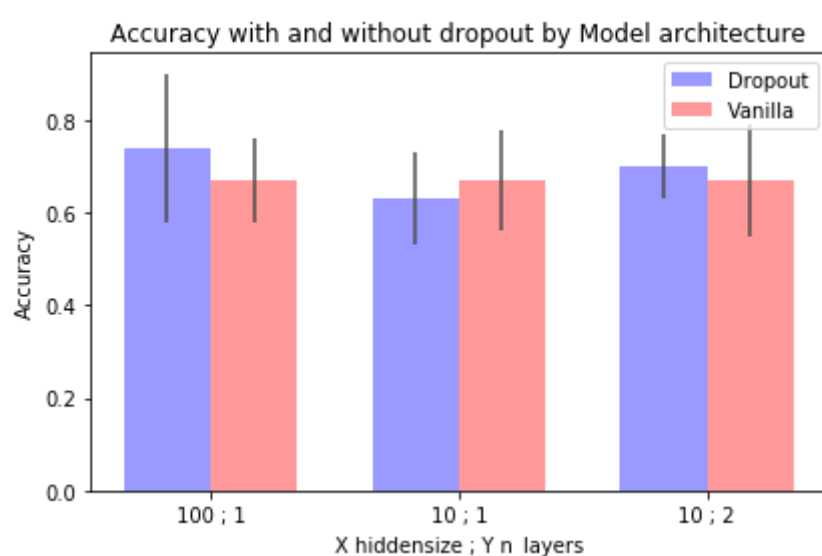


**Fig. Continuous learning. Left : no, right : yes.**

### 3 Regularization

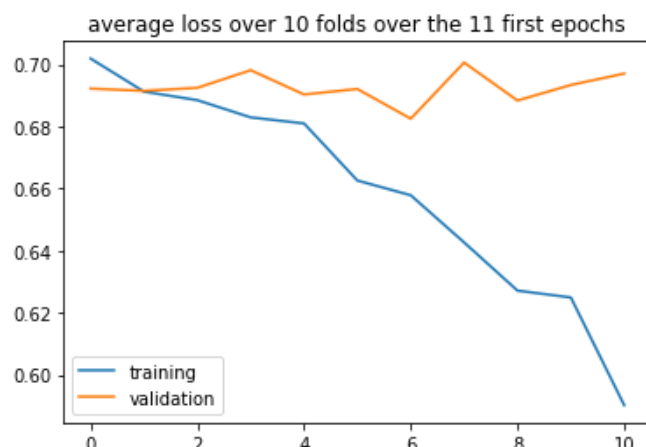
Consistently with Lipton et al., adding a dropout of 0.5 between the non-recurrent layers enhances the results with large layers and stacked LSTMs. See figure below, all the other hyperparameters are fixed to :

dowsampling factor	learning_rate	bidirectional	carry over	gradient clipping
1	0.001	TRUE	0	5





Adding a dropout of 0.5 between the non-recurrent layers helped with overfitting, with and without downsampling.



Interestingly, this is specific to LSTM, using feedforward dropout on GRU had very little effect...

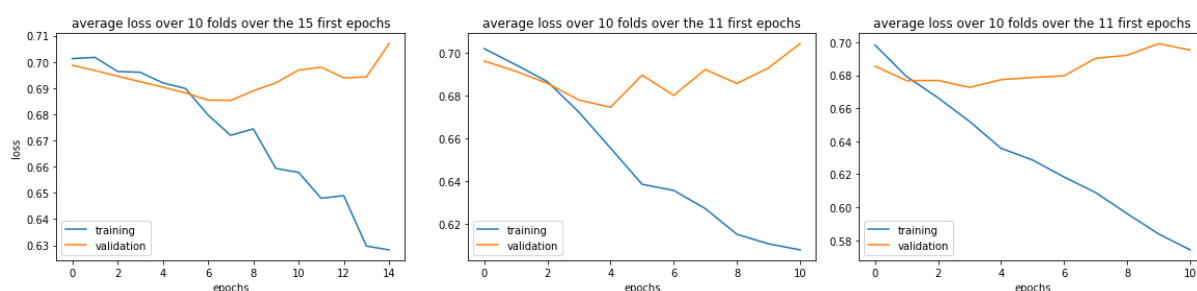
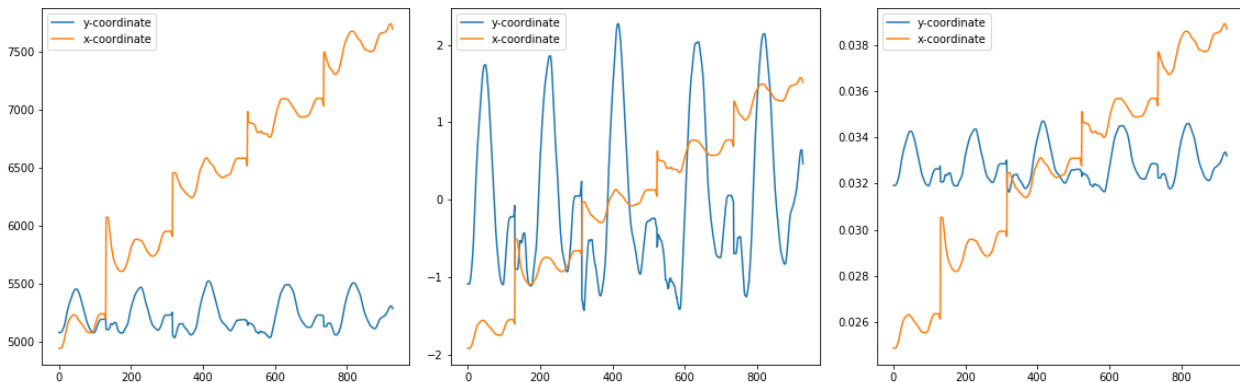


Fig. From left to right : 1GRU with and without dropout, bGRU without dropout (different scales)

## 4 Data preprocessing and representation

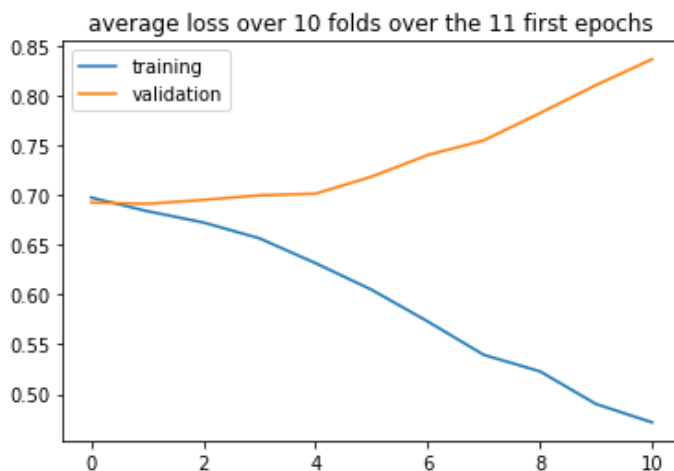
Standardizing a vector (e.g. x-coordinate of a task) gives it a mean of 0 and a standard deviation of 1. Normalizing a vector gives it a norm of 1.



**Figure of the x and y coordinate recorded at 200Hz for the second task (writing multiple "I") of the first subject (PD) : raw, standardized and normalized (from left to right).**

One can notice that the data ranges are preserved when normalizing unlike when standardizing. So I'd assume that it was better to normalize than to standardize but after experiment, the model won't fit the data if it's normalized, although it fits it quite well if it's standardized. After investigation I noticed that, between subjects, the normalization shifted the values of the data : you can see here that for the first subject the normalized x coordinate ranges between 0.026 and 0.038, whereas for the 26th subject, it ranges from 0.022 to 0.028 (although the raw ranges are similar). If we feed the raw data to the model (i.e. without any type of normalisation) it doesn't fit to the data : with  $lr = 10^{-2}$ , the loss converges at 0.7 (accuracy 0.469). Tried with  $lr = 10^{-3}$  (converges higher),  $lr = 10^{-1}$  (doesn't converge).

Downsampling the data 10 times doesn't help with the overfitting as you can see on line 14.



The model is *not* able to fit the data when concatenating a one hot vector to the measures which identifies the task (e.g. [0., 1., 0., 0., 0., 0., 0., 0.] for the *l* task). This might be

considered an early fusion technique which augments the dataset. We might want to compare it to late fusion techniques (e.g. majority voting).

In order to be able to implement *attention* (cf. Report #4 [1.1 Model Architectures et al.](#)) I tried trimming and padding with zeros the spirals that were respectively longer and shorter than 3117 timesteps (3rd quartile of the spirals' length) and the model won't fit the data with the same hyperparameters as without trimming and padding.

Inspired by the work of Zhang et al. (cf. report #5). I tried to split the model into subsequence with a fixed time window of 100 timesteps. Unfortunately, the model is completely unable to generalize and the validation accuracy never gets better than chance level, although the model fits the training set. Dropout was ineffective. I should maybe try other segmentations, e.g. stroke, also I might want to train 2 separate models on in-air and on-paper strokes, respectively. This is motivated by the experiment results which are superior on the *spiral* task where the single stroke is on-paper. As opposed to the *l* task where strokes are both on paper and in-air. Since in-air movements are very different from on-paper ones, I think they might disturb the model.

Interestingly, computing the movement (cf. report #5) had no effect when the data was split in subsequences but improved the results when training on the whole sequence, although Leo told me it shouldn't be necessary when using a RNN, therefore, further experiments are required.

## 5 Evaluation and experiments

Average metric over the 10 folds (+ standard deviation) after early stopping (i.e. we select the best epoch based on the validation accuracy for each fold). For the baseline model depicted in [1 Model architecture](#).

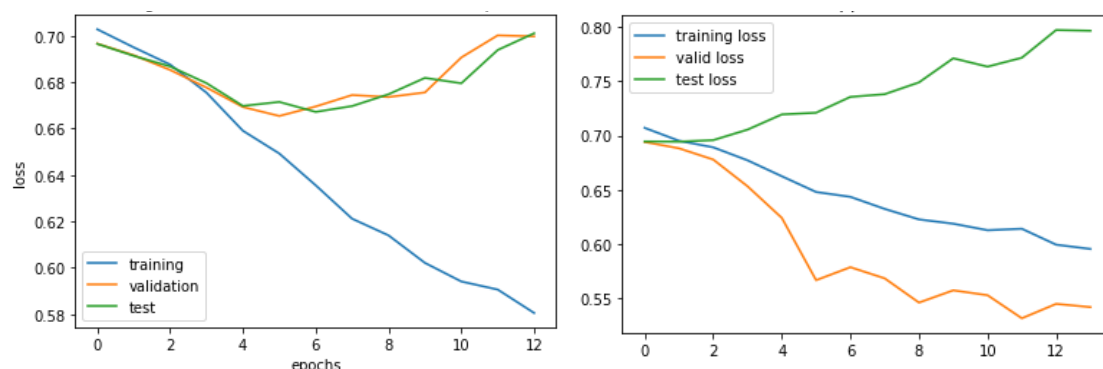
TRAIN accuracy	accuracy	Se	Sp	PPV	NPV
0.59 (+ 0.07)	0.68 (+ 0.09)	0.54 (+ 0.20)	0.82 (+ 0.20)	0.81 (+ 0.20)	0.65 (+ 0.08)

Every metric is for the validation set if not specified otherwise (i.e. except for the first column). Since the validation set is very small (~ 7 subjects) the metric standard deviation (std) is very high, this is consistent with Moetesum et al. (cf. Report #2). The accuracy for the *spiral* task is 76%, 63% and 55% for Moetesum et al., Drotar et al. and Impedovo et al., respectively. Notice that, in the same way as Drotar et al., we validate on the test set (i.e. on the validation fold). Therefore our results are over optimistic and possibly overfitted to the dataset.

I tried to train the model without early stopping for 4 epochs (because the model usually fits after 4 epochs) with the GRU depicted in [1 Model architecture](#). The results are slightly decreased and the std augments : from 70% to 65% and from 8% to 13% respectively. Moreover, these results are also "overfitted" to the test set because if we train for one more epoch they get worse so I didn't present them.

I also tried to implement a proper early stopping by splitting the dataset in training, validation and test set and early stopping according to the validation results. However in this case the results are even worse, only 61% accuracy and 17% std.

I don't think the problem is that the training set is too small because if we take the best epoch for each test folds we get the same results as before : 71% (+ 11%) accuracy. The problem is that since the validation test is very small, it's not very representative of the generalization capacities of the model. On average the validation loss and the test loss are quite similar (cf. figure below) but on some folds they're very different. We might try to do data augmentation on the validation set so it's bigger and thus more representative of the capacities of the model.



**Fig. Left : average loss over the 10 folds, right : loss on the first fold.**

Even though these results are disappointing, I think we should stick with one of these methods because Drotar et al. validate on the test set for hyperparameters search (as us), not for the  $n^\circ$  of epochs etc. I hope we will be able to conjecture from the previous experiments results if we keep on early stopping since the dataset split is now different. If the data augmentation does not improve the results I think we should not perform early stopping.

Following the guidelines of Greff et al. (cf. Report #4 [1.2 Hyperparameters finetuning](#)), we should be able to fine tune the hyperparameters independently, starting from :

- learning rate
- hidden size
- number of layers

- regularization techniques (e.g. dropout, although Karpathy<sup>2</sup> and Lipton et al. suggests that it interacts with the hidden size so we might want to tune these hyperparameters together)

Throughout these experiments we shall use a bidirectional LSTM, which can only better the results in my opinion (cf. [1 Model architecture](#)). Also the decoder architecture will be kept as simple as possible (cf. [1 Model architecture](#)). I made this choice because I assume that the "hard part of the job" is made by the LSTM and not the decoder.

I thought we might also use an SVM as Moetesum et al. or OPF as Passos et al. 2018 (cf. Report #2) but we can't because unlike them we need to train the network (they use a CNN pretrained on ImageNet).

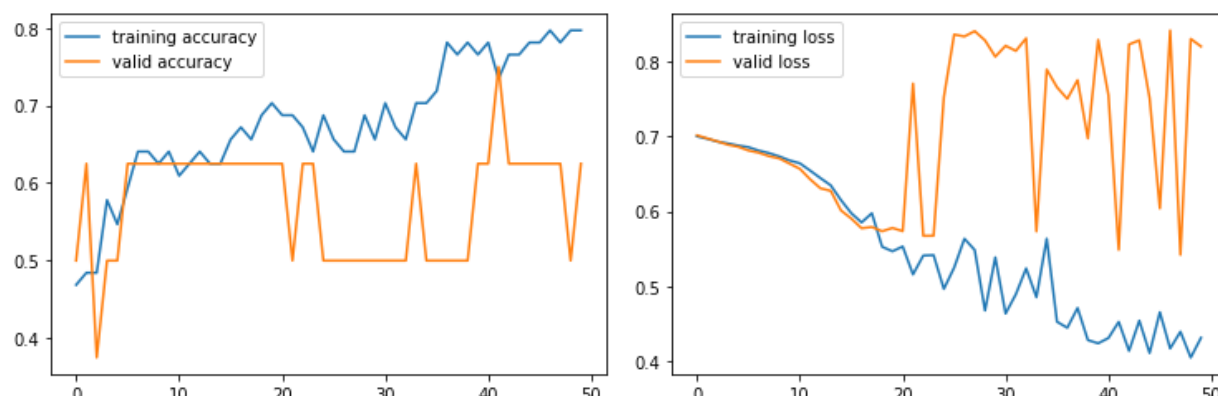
All the results of the experiments are available in the git repo. The studied hyperparameter is print in **bold**.

Note :

- the max epochs was set at 5 for the first 2 lines, 10 for the 3rd and 4rd lines and 50 for the rest
- there was no patience for early stopping for the lines 2 to 8 but no improvement was okay : you can see the differences between line 8 and 9. They are probably explained because when we tolerated the lack of improvement, the model could wait for tens of epochs (cf. "early stopped"), overfitting until "by chance" the accuracy would went up (cf. figure below : left for accuracy, right for loss). Thus the **results of lines 8 are overoptimistic**.
- the seed for random weights initialization was not fixed for the lines 2 to 8. Therefore a direct comparison with and between those might be biased.
- on the lines 2 to 12, all the biases were init randomly. You can see the difference on line 13 where we started to init the forget gate bias at 1 : it improved training accuracy, validation PPV and reduced the variance (std).
- Some "falses" are greyed out, it's because they represent subject indexes before discarding the 3 subjects who didn't perform the spiral task. Therefore these index don't match the other ones.

---

<sup>2</sup> <https://github.com/karpathy/char-rnn>



**Fig. Model overfitting until (probably) random amelioration (line 8). Left : accuracy, right : loss.**

## 6 Visualization & Interpretation

### 6.1 Input weights

Karpathy, Johnson et al. wrote a great article about visualizing RNN, however, their classification task (language modeling) allows for a simpler interpretation. I tried visualizing the input weights of a model fitted to the data in order to see which of the 7 measures was activated. As you can see below, the weights differ along the layer (i.e. between the LSTM units) without any measure standing out (I also tried plotting each measure along the layer or compute the mean of each measure). One possible interpretation is that there's no measure which is consistently "interesting". However it's quite strange given that this model was trained only on the *spiral* task where the button status (column n°3 here) is always equal to 1.

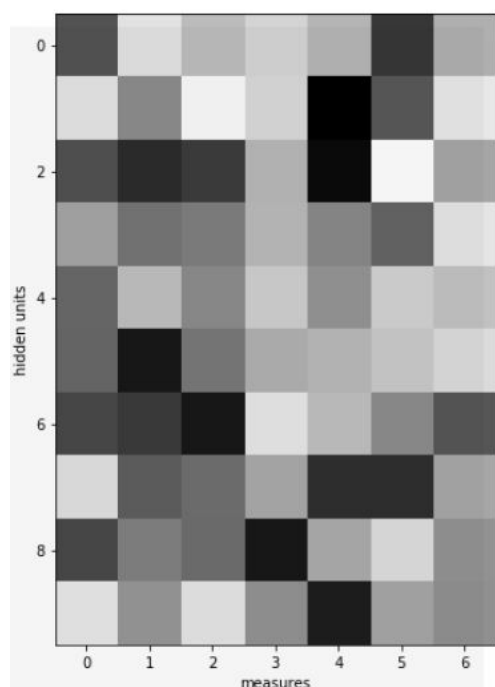
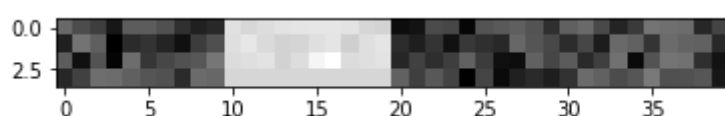


Fig. Plot of the weights of the 10 first input units (of the one and only layer of the network).

## 6.2 Forget Gate

Interestingly, after using the trick of init the forget gate bias at 1, even after 22 epochs and an overfitted model, the bias of the forget gate are still worth  $\sim 1$  (cf. figure below). One possible interpretation is that the model learned to *remember* the previous value of the *cell state*.



## 6.3 Misclassified subjects

For the model with hyperparameters :

dowsam pling factor	learning_ rate	hidden_s ize	num_lay ers	bidirectio nal	carry over	dropout	gradient clipping
1	0.001	50	1	TRUE	0	0.5	5

I looked into which subjects had been misclassified when the model had early stopped (i.e. when the classification accuracy was highest). It turns out that, among the misclassified subjects (22 over all 10 folds), 3 of them had performed very long exams of 16071, 9724, and 8581 timesteps. Thus bringing the average length of misclassified subjects to 3552 instead of 2758 (average over every subjects). It'd be interesting to see if the same subjects are misclassified across models.

## Conclusion - Todo List

Cf. report #4 conclusion and todo list.

I followed my intuition and explored the hyperparameters fine-tuning of our current model on the *spiral* task. I found that GRUs were less prone to overfitting than LSTMs, although feedforward dropout was efficient for regularizing LSTM. Thus, stacked LSTMs give better results than long, single layers.

In the experiment results, I colored in green the hyperparameters that worked well together, as well as all the learning rates and gradient clipping which I found the best values to be  $10^{-3}$  and 5, respectively. Again, all the experiments were conducted while validating on the test set (cf. [5 Evaluation and experiments](#)) so I hope the results are still valid when running a proper CV.

Then I tried, without success, to learn from all tasks at the same time, then to learn from fixed-time subsequences (cf. [4 Data preprocessing and representation](#)). My next move a priori is data augmentation (cf. report #4).